

Retrotransposon mobilization in cancer genomes

Tracy Ballinger¹, Adam D. Ewing¹, David Haussler^{1,2}

¹Center for Biomolecular Science and Engineering, University of California, Santa Cruz, CA 95064, USA.

²Howard Hughes Medical Institute, University of California, Santa Cruz, CA 95064, USA.

January 20, 2015

Abstract

The Cancer Genome Atlas project was initiated by the National Cancer Institute in order to characterize the genomes of hundreds of tumors of various cancer types. While much effort has been put into detecting somatic genomic variation in these data, somatic structural variation induced by the activity of transposable element insertions has not been reported. Transposable elements (TEs) are particularly relevant in cancer in part because of several known cases in which a TE insertion is directly linked to cancer formation and studies linking the epigenetic status of retrotransposons to carcinogenesis and patient outcome. Additionally, evidence for somatic retrotransposition in eukaryotic genomes suggests that some tissues and therefore some cancer types may be disposed to increased retrotransposition. We built upon previous work to develop a highly efficient computational pipeline for the detection of non-reference mobile element insertions from high-throughput paired-end whole genome sequencing data that is capable of detecting breakpoints through a local assembly strategy. Using this, we analyzed 33 whole genome tumor datasets with paired normal samples from TCGA across 3 different cancer types: glioblastoma multiforme (GBM), ovarian serous cystadenocarcinoma (OV) and colorectal adenocarcinoma (COAD). We detected 72 insertions in colon samples, almost all of them LINE-1 elements, and none in GBM or OV. The amount of somatic retrotransposition varies widely between samples with 61 insertions present in one case. The lack of somatic retrotransposon insertions in GBM and OV samples suggests that TE activity in cancer is restricted to certain cancer types.

1 Introduction

Retrotransposons are found in all eukaryotic genomes. They are observed as repetitive DNA elements due to their capacity to insert new copies of themselves into the host DNA through a copy and paste process using an RNA intermediate (Boeke et al., 1985). They are categorized as either long terminal repeat (LTR) or non-LTR and further into families based on sequence similarity to other elements and by their mechanism of mobilization. The non-LTR retrotransposons that inhabit mammalian genomes are likely to mobilize through a mechanism known as target-primed reverse transcription (Luan et al., 1993). Numerous retrotransposon copies exist in the human genome, comprising at least 45% of its DNA (Lander et al., 2001) and perhaps over two-thirds when highly sensitive TE detection methods are applied (de Koning et al., 2011). The most prolific retroelements in the human genome include LINE-1 and Alu sequences, comprising 17% and 8% of the assembled reference genome, respectively. During primate evolution, the general pattern of retroelement activity has been for one family of LINE element to be active at a time, suggesting either competition for a host factor or adaptation to evade one (Khan et al., 2006). LINE-1 elements are autonomous retrotransposons encoding two proteins (Scott et al., 1987) responsible for both

their own mobilization *in cis* and the mobilization of non-autonomous Alu elements (Dewannieux et al., 2003), SVA elements (Hancks et al., 2011), and processed pseudogenes (Esnault et al., 2000) *in trans*. The activity of the human-specific LINE element, termed L1HS, was first recognized *in vivo* due to its ability to disrupt exons and cause Mendelian disease (Kazazian et al., 1988). Since then, transposable elements have been linked to a variety of diseases, including cancer, through insertional mutagenesis of exons and regulatory regions near genes, disrupting gene function or regulation (see Hancks et al. 2012 for review). For example, in one case, an exonic L1 insertion was found in the APC tumor suppressor gene in colon cancer tissue but not the normal tissue of the same patient (Miki, 1992). Intronic retroelement insertions are known to affect splicing by providing 5' or 3' splice sites or disrupting sequence at the branch point (Belancio et al., 2008; Hancks et al., 2009; Taniguchi-Ikeda et al., 2011). Recent estimates place the rate of L1 retrotransposition in human genomes at 1 new insertion for every 100 to 150 live births (Ewing and Kazazian, 2010; Huang et al., 2010). Since retroelements can clearly have an impact on phenotype and disease, it is important that retrotransposon insertion polymorphisms (RIPs) and mutations be characterized in genomic studies. A plethora of recent studies provide various means to document retrotransposon insertion polymorphisms (RIPs) segregating in human populations (Beck et al., 2010; Ewing and Kazazian, 2010; Hormozdiari et al., 2010a; Huang et al., 2010; Iskow et al., 2010; Witherspoon et al., 2010; Ewing and Kazazian, 2011; Stewart et al., 2011), including one report of 9 somatic retrotransposon insertions across 6 lung tumors (Iskow et al., 2010).

Cancer progression depends on the accumulation of somatic mutations, and recent evidence suggests that retrotransposition also occurs in some somatic tissues such as neuronal stem cells (Muotri et al., 2005; Coufal et al., 2009; Baillie et al., 2011). The observation of somatic retrotransposition in specific tissue types suggests tissue-specific regulation, either through known regulators such as APOBEC3 proteins (Kinomoto et al., 2007; Muckenfuss et al., 2006; Stenglein and Harris, 2006; Chen et al., 2006), germline piRNAs (Aravin et al., 2007), and DNA methylation (Yoder et al., 1997; Bourc'his and Bestor, 2004), or through novel mechanism(s) not yet ascribed to transposable elements. Other lines of evidence for somatic retrotransposition include the aforementioned disease-causing insertions, observations of varying levels of transgene-borne somatic retrotransposition in transgenic mice (Kano et al., 2009), and somatic R2 insertions in *Drosophila simulans* (Eickbush and Eickbush, 2011). In addition to mutagenizing both somatic and germline genomes through new insertions, transposable elements play an important role in shaping gene regulatory networks by providing binding sites for transcription factors, including those highly important for cancer progression such as *TP53* and *SOX2* (Wang et al., 2007; Bourque et al., 2008; Kuwabara et al., 2009; Harris et al., 2009). Furthermore, genome-wide methylation status is often assessed through analysis of CpG islands located in the 5' UTRs of LINE-1 elements, which are typically heavily methylated (Woodcock et al., 1997), contributing to their quiescence in most somatic tissue types. Through this assay, a wide variety of cancers are found to be hypomethylated (Ogino et al., 2011), leading us to speculate that retrotransposition rates may be substantially increased in certain cancer types or samples.

In order to test this hypothesis, we took advantage of whole-genome sequence data available through The Cancer Genome Atlas (TCGA). TCGA is an ongoing multi-institutional effort that will eventually include whole genome sequence data for hundreds of tumors and corresponding normal samples for over 20 different cancer types. Here, we consider transposable element insertions in the genomes of three cancer types: glioblastoma multiforme (GBM) (TCGA Research Network, 2008), ovarian serous cystadenocarcinoma (OV) (TCGA Research Network, 2011), and colon/rectal adenocarcinoma (COAD/READ) (Muzny et al., 2012), and present evidence for substantially increased retrotransposition in colorectal adenocarcinoma.

2 Results

We developed a computational pipeline (discord-retro, <http://github.com/adamewing/discord-retro>) to detect non-reference retrotransposon insertions from paired end whole genome sequencing data by using mate-pair information from discordantly mapped read pairs (see Methods). We measured the detection characteristics of our application by repeatedly inserting 100 retrotransposons into the euchromatic portion of human chr22 at random positions, generating paired reads, mapping to the GRCh37 reference sequence, and applying our method (see Methods for details). We observe 87.9% sensitivity with perfect specificity when insertions into other insertions of the same class (e.g. LINE into a LINE or Alu into an Alu) are discarded, and 94.5% sensitivity and perfect specificity if these insertions are allowed. Using discord-retro, we analyzed 33 high coverage ($>30\times$) tumor and normal genome pairs produced by TCGA, and identified retrotransposon insertions not found in the human reference genome (NCBI36 or GRCh37). For high coverage data, the tumor and patient-matched normal paired-end sequence data were combined in order to distinguish between a non-reference germline insertion, which would be found in both tissues, from a somatic insertion, which would be found only in the tumor (or normal) DNA. Refinement of junctions using local assembly and analysis of soft-clipped reads allowed breakpoint-resolution on one or both ends of a predicted insertion. In 45% of cases, we identified target site duplications (TSDs), a short duplication of sequence around the insertion site that occur as a byproduct of target primed reverse transcription.

2.1 Germline Insertions

Across the 33 tumor/normal pairs, we identified 7022 non-reference retrotransposon insertions present in both the tumor and corresponding normal genome, of which 3273 overlapped with a previous study (Wang et al., 2006; Beck et al., 2010; Ewing and Kazazian, 2010; Hormozdiari et al., 2010a; Iskow et al., 2010; Witherspoon et al., 2010; Ewing and Kazazian, 2011; Stewart et al., 2011) (Fig S1). Of all insertions detected, 727 were LINE insertions, 6101 were Alu insertions and 189 were SVA insertions. Of the 3749 previously uncataloged insertions, 350 were LINE-1 elements, 3220 were Alu elements, and 177 were SVA elements. For every tumor/normal pair we detected an average of 111 LINEs, 823 Alus, and 36 SVAs (Fig. 1). We detected an average of 8.4 LINE, 79.6 Alu and 1.9 SVA insertions that were present in only one sample. The chromosomal distribution of insertions is illustrated in Figure 2.

2.2 Somatic Insertions

Somatic insertions are those occurring exclusively in either the tumor or normal sample of a patient-matched pair of genomes and also not present in any other sample or catalogue of retrotransposon insertions from a previous study. Furthermore, because we combine discordant read pairs across both the tumor and normal tissue for an individual, we can be sure that if a tumor-specific or normal-specific call is made, not a single read that could indicate the presence of the insertion exists in the other sample. We found 72 tumor-specific LINE-1 insertions from 4 colon cancer tumor/normal pairs and none from any of the 18 GBM and 10 OV tumor/normal pairs (Fig. 3). Conversely, for the four samples with tumor-specific insertions, we found very few insertions present only in the normal sample (Fig. S2). For one sample, TCGA-AA-A00R, there was an abnormally high number of normal-specific predictions which we believe to be an artifact. The distribution of insertion lengths for the tumor-specific LINE-1 insertions differs markedly from insertions found in both the tumor and normal tissues (Fig. 4). Of 655 LINE-1 insertions present in both tumor

and normal genomes, 147 were full length, defined here as 5.8 kb or greater as indicated by the minimum and maximum mapping locations of paired reads within the reference elements, which resembles the distribution expected from the length distribution in the reference genome (Grimaldi et al., 1984; Pavlíček et al., 2002). In contrast, only 2 out of 72 tumor-specific LINE-1 insertions were full-length ($p < 9.68 \times 10^{-6}$, Fisher’s exact test).

Of the 72 tumor-specific insertions, one occurred in the 3’ UTR of *PPP1R1C* (protein phosphatase 1, regulator subunit 1C), and 22 occurred in introns, including an insertion in *NAV3*, a gene associated with colon cancer (Carlsson et al., 2012). The *NAV3* insertion occurs in patient TCGA-AA-3518 115 bp downstream of the third exon in the same orientation as the gene in a region overlapped by DNase hypersensitivity and H3K27 acetylation signals (Fig. S3), indicating possible regulatory elements nearby. An examination of gene expression levels (Agilent 244K Custom Gene Expression G4502A-07-3) between the tumor and normal insertions was carried out using the UCSC Cancer Genome Browser (Zhu et al., 2009), which indicated lower expression of *NAV3* and other genes containing tumor-specific insertions (notably *A2BP1*, and *CTNNA2*) in the tumor relative to the patient-matched normal colon tissue (Fig. 5). Here we focus further analysis on *NAV3*, as decreases in its expression are frequent in colorectal adenomas (Carlsson et al., 2012). Analysis of Agilent expression data from TCGA-AA-3518 shows the difference in expression between *NAV3* in tumor and normal tissues is ranked at the 93rd percentile relative to differences between all other genes in the same tumor/normal pair. To ascertain whether other somatic mutations might be responsible for the observed change in expression, we compared *NAV3* and the surrounding region between cancer and normal genomes of TCGA-AA-3518 using the BamBam algorithm (Z. Sanborn, unpublished). We detected evidence for 9 potentially cancer-specific SNPs in the 382 kbp region (Fig. S4), but we found no evidence for point mutations or CNVs in exons or in the proximal upstream region likely to have an obvious effect on transcript abundance apart from the L1 insertion.

2.3 Similarity to reference elements

After acquiring a set of insertion predictions for each sample, we sought to determine the closest element in the reference genome in terms of sequence similarity, as this may represent an element similar to the active progenitor element. In general, it is unlikely that the true progenitor can be identified through sequence similarity alone, as the active elements in the human reference genome diverge from one another by less than 1% (Brouha et al., 2003; Seleme et al., 2006; Beck et al., 2010). That said, identification of the most similar elements in the human reference genome based on local re-assembly of the elements detected by discord-retro yields an enrichment of full-length, intact, human L1 elements, some of which are known to be active elements (Table 1). This serves as further evidence to substantiate our claim that our report of 72 tumor-specific L1 insertions in 4 colon cancer cases are novel insertions derived from active L1 elements.

3 Discussion

As TCGA and others continue to sequence more cancer and paired normal cases across a wider number of cancer types, we may uncover clear driver mutations caused by transposable elements and other cancer types that exhibit high levels of insertional mutagenesis by transposable elements. It is remarkable that TE activity appears so much higher in colorectal adenocarcinomas than in glioblastoma multiforme or ovarian serous cystadenocarcinoma, but the specific mechanism behind this tissue specificity has eluded us so far. An intronic insertion in *NAV3* seems to be paired with a marked decrease in gene expression, although it is far from clear if there is a direct relationship

between the presence of somatic LINE-1 insertions and the expression decrease in cases like this. Given that a survey of other somatic mutations in and surrounding *NAV3* yielded nothing that stood out as a possibly expression-altering mutation, and an insertion in *NAV3* occurred only 112 bp downstream of an exon in a region with an epigenetic profile indicating regulatory potential, we posit that the LINE-1 insertion may be responsible for the cancer-specific decrease in expression in this instance.

Our general knowledge of somatic retrotransposition is limited, although new technologies and the decreasing cost of sequencing will likely provide new insights in the near future as sequencing studies begin to focus on multiple tissues from a single donor individual. In most respects, the somatic tumor-specific insertions detected by our method are similar to germline insertions, with the notable exception of their length: 97.2% of tumor-specific insertions are truncated as compared to a 77.6% truncation rate for germline insertions (insertions detected in both normal and cancer samples). The mechanism for L1 truncation is unknown; conjectures include the presence of an endo- or exonuclease that targets the L1 RNA template, or a factor that interferes with reverse-transcription. At this stage the etiology of element truncation in colon tumors and how it differs from normal or germline tissue is unknown. It may be illuminating to work out why tumor-specific insertions are more severely truncated than those in the germline as a future study and whether this is a general characteristic of somatic retrotransposition or if there is some connection to tumor biology.

This is an exciting time for transposable element biology given our improving ability to explore entire genomes. In this case, whole-genome paired-end sequencing has allowed us to detect somatic retrotransposition in cancer genomes, an observation that opens many new questions regarding the role of mobile DNA in carcinogenesis and tumor molecular biology. As sequencing technologies and our ability to detect structural variants improve, so will our ability to characterize new TE insertions and their parent elements, perhaps gaining further insight on what leads to tissue or disease specific TE activation.

4 Methods

A number of successful computational methods have been devised capable of detecting transposable element insertions from whole-genome sequence data including VariationHunter2 (Hormozdizari et al., 2010b), T-lex (Fiston-Lavier et al., 2011), RetroSeq (<https://github.com/tk2/RetroSeq>), HYDRA-SV (Quinlan et al., 2010), and Tea (Lee et al., 2012). The approach outlined here, implemented as discord-retro, has the advantage of working directly from the ubiquitous .bam sequence alignment format (as does RetroSeq) with minimal need for additional mapping apart from that required to identify insertion breakpoints. Here, we give a high-level overview of our method. Sequence data analyzed in this study was generated on the Illumina platform and aligned to a human reference assembly (NCBI36 or GRCh37) by TCGA Research Network members at TCGA Genome Sequencing Centers.

Paired-end reads can be classified based on how they map to the reference genome. A read pair is called concordant if both reads map the proper distance apart and in the correct orientation for the insert size and procedure used in the library preparation and sequencing, and discordant if these conditions are not met. For example, ends of a discordant paired read may map to different chromosomes, too far apart, too close together, or in the wrong orientation. A second type of improperly paired reads are ones in which one read maps to the reference, but its pair does not. These reads are referred to as one-end-anchored (OEA). Lastly, reads are called soft-clipped when part of the read aligns to the reference sequence, but either or both ends of the read do not.

We first selected all discordant reads from both the tumor and normal sequencing data of a patient where one read of the pair maps to a unique portion of the genome, called the “anchored” read, and the other end maps to a repeatmasker annotation elsewhere in the genome. We will refer to these types of read-pairs as one-end-repeat (OER) reads. We filtered elements corresponding to AluS and LTR elements from the results due to an overabundance of calls with no corresponding breakpoint predictions in some samples. Regions where the uniquely mapped ends of the OER reads clustered in two peaks with opposite orientation were considered consistent with an insertion existing between the two clusters of OER reads. We require there to be 8 OER read pairs within a 500bp window, and for there to be at least 2 uniquely mapped or “anchored” reads on either strand. The requirement that both breakpoints (5’ and 3’ junctions) be covered by paired reads reduces the chance of incorrectly annotating a segmental duplication, translocation, or inversion as a transposable element insertion.

The selection of clustered discordant OER reads yields a set of 20-50bp windows as predicted transposable element insertion sites. These were annotated as “germline” if there were discordant reads in both the tumor and normal tissue samples, as “somatic/cancer” if there were contributing discordant reads only in the tumor tissue, and as “somatic/normal” if there were contributing discordant reads only in the normal tissue. Insertion loci are cross-referenced against retrotransposon insertion polymorphisms (RIPs) cataloged in previous studies (Wang et al., 2006; Beck et al., 2010; Ewing and Kazazian, 2010; Hormozdiari et al., 2010a; Iskow et al., 2010; Witherspoon et al., 2010; Ewing and Kazazian, 2011; Stewart et al., 2011) and against each other. As breakpoint resolution varies across studies, insertions within the same 500bp window were considered overlapping. A total of 14 (16%) potentially tumor-specific insertions were eliminated by comparison to known RIPs and RIPs found in this study.

4.1 Breakpoint refinement using soft-clipped reads

Soft-clipped reads mapped using bwa (Li and Durbin, 2009) could be used to pinpoint a breakpoint in the insertion site. For each of the 14 samples aligned with bwa (Table S1), soft-clipped reads mapping within 500bp of each of the predicted insertion sites and which had greater than 10 bp clipped from the read were used to find a consensus breakpoint where a majority of soft-clipped read ends occurred at the same nucleotide in the reference genome. When breakpoints for both the 5’ and 3’ junctions between the element and the reference genome were detected, we identified target site duplications when the breakpoint on the forward strand occurred 3-50bp downstream of the breakpoint on the reverse strand.

4.2 Breakpoint refinement with local assembly

We used a local assembly and realignment strategy to determine breakpoints for all samples. All discordant and soft-clipped reads within 500bp of a predicted insertion site were assembled using Velvet (Zerbino and Birney, 2008) with a k-mer size of 31, the shortPaired option, and insert length of 300. If the reads assembled into 5 contigs or less, these contigs were mapped back to the reference genome using BLAT. A cutoff of 5 contigs was chosen because when more contigs were present, they were generally too short to be more informative than the original reads. After mapping the assembled contigs back to the reference assembly, breakpoints present as the point where a contig no longer matches the reference sequence and begins matching a reference retroelement sequence. Target site duplications could be ascertained in cases where two assembled contigs had overlapping alignments to the predicted insertion site on opposite strands.

4.3 Simulation

To measure the accuracy and sensitivity of our pipeline, we inserted 100 LINE, SINE, and SVA sequences randomly into the euchromatic sequence of chr22 from hg19/GRCh37. The retroelement sequences were randomly truncated on the 5' end up to 75% of the original element length for LINEs and SVAs, and 25% for Alus. Poly(A) tails between 20 and 70bp in length were added to the 3' end, and 12bp of the target insertion site was duplicated on the 5' junction to mimic target site duplications. Paired Illumina reads were simulated via wgsim (<https://github.com/lh3/wgsim>) to generate paired 75bp reads at 30x coverage. These reads were mapped back to the reference genome using bwa (Li and Durbin, 2009) with the following parameters: `-q 5 -l 32 -k 2 -t 4 -o 1`, alignments were processed with samtools (Li et al., 2009) and used as input to discord-retro.

4.4 Assessing sequence similarity to reference elements

We performed local sequence assembly as described in section 4.2 to generate contigs corresponding to inserted sequences. BLAT alignments of the contigs were carried out to find the most closely related elements in the reference genome. Repeatmasker-annotated elements were scored by the sum of the products of the percent identity of a BLAT alignment and the length of the alignment for each contig that overlapped the repeat masked element. The element with the highest score was predicted to be the source element for the new insertion, excluding elements within 1000bp of the insertion site. Elements scoring within 20 percent-identity bases of the highest score were considered as potential progenitors as well. In cases where there were several repeat elements tied for the highest score or very close to the highest score, the progenitor is considered ambiguous. We ranked the repeat masked elements by the number of times they were predicted to be a progenitor for an somatic insertion, whether ambiguous or not, and examined the top 10 elements for retrotransposition capability.

References

- Aravin, A. A., Hannon, G. J., and Brennecke, J., 2007. The Piwi-piRNA pathway provides an adaptive defense in the transposon arms race. *Science (New York, N.Y.)*, **318**(5851):761–4.
- Baillie, J. K., Barnett, M. W., Upton, K. R., Gerhardt, D. J., Richmond, T. A., Sapio, F. D., Brennan, P., Rizzu, P., Smith, S., Fell, M., *et al.*, 2011. Somatic retrotransposition alters the genetic landscape of the human brain. *Nature*, .
- Beck, C. R., Collier, P., Macfarlane, C., Malig, M., Kidd, J. M., Eichler, E. E., Badge, R. M., and Moran, J. V., 2010. Line-1 retrotransposition activity in human genomes. *Cell*, **141**(7):1159–1170.
- Belancio, V. P., Roy-Engel, A. M., and Deininger, P., 2008. The impact of multiple splice sites in human L1 elements. *Gene*, **411**(1-2):38–45.
- Boeke, J. D., Garfinkel, D. J., Styles, C. A., and Fink, G. R., 1985. Ty elements transpose through an rna intermediate. *Cell*, **40**(3):491–500.
- Bourc'his, D. and Bestor, T. H., 2004. Meiotic catastrophe and retrotransposon reactivation in male germ cells lacking dnmt3l. *Nature*, **431**(7004):96–9.
- Bourque, G., Leong, B., Vega, V. B., Chen, X., Lee, Y. L., Srinivasan, K. G., Chew, J.-L., Ruan, Y., Wei, C.-L., Ng, H. H., *et al.*, 2008. Evolution of the mammalian transcription factor binding repertoire via transposable elements. *Genome Res*, **18**(11):1752–62.

- Brouha, B., Schustak, J., Badge, R. M., Lutz-Prigge, S., Farley, A. H., Moran, J. V., and Kazazian, H. H., 2003. Hot L1s account for the bulk of retrotransposition in the human population. *Proceedings of the National Academy of Sciences of the United States of America*, **100**(9):5280–5.
- Carlsson, E., Ranki, A., Sipilä, L., Karenko, L., Abdel-Rahman, W. M., Ovaska, K., Siggberg, L., Aapola, U., Ässämäki, R., Häyry, V., *et al.*, 2012. Potential role of a navigator gene NAV3 in colorectal cancer. *British journal of cancer*, **106**(3):517–24.
- Chen, H., Lilley, C. E., Yu, Q., Lee, D. V., Chou, J., Narvaiza, I. n., Landau, N. R., and Weitzman, M. D., 2006. APOBEC3A is a potent inhibitor of adeno-associated virus and retrotransposons. *Current biology : CB*, **16**(5):480–5.
- Coufal, N. G., Garcia-Perez, J. L., Peng, G. E., Yeo, G. W., Mu, Y., Lovci, M. T., Morell, M., O’Shea, K. S., Moran, J. V., and Gage, F. H., *et al.*, 2009. L1 retrotransposition in human neural progenitor cells. *Nature*, **460**(7259):1127–31.
- de Koning, A. P. J., Gu, W., Castoe, T. A., Batzer, M. A., and Pollock, D. D., 2011. Repetitive elements may comprise over two-thirds of the human genome. *PLoS genetics*, **7**(12):e1002384.
- Dewannieux, M., Esnault, C., and Heidmann, T., 2003. LINE-mediated retrotransposition of marked Alu sequences. *Nature genetics*, **35**(1):41–8.
- Eickbush, M. T. and Eickbush, T. H., 2011. Retrotransposition of R2 elements in somatic nuclei during the early development of *Drosophila*. *Mobile DNA*, **2**(1):11.
- Esnault, C., Maestre, J., and Heidmann, T., 2000. Human LINE retrotransposons generate processed pseudogenes. *Nature genetics*, **24**(4):363–7.
- Ewing, A. D. and Kazazian, H. H., 2010. High-throughput sequencing reveals extensive variation in human-specific l1 content in individual human genomes. *Genome Res*, **20**(9):1262–1270.
- Ewing, A. D. and Kazazian, H. H., 2011. Whole-genome resequencing allows detection of many rare LINE-1 insertion alleles in humans. *Genome research*, **21**(6):985–90.
- Fiston-Lavier, A.-S., Carrigan, M., Petrov, D. A., and González, J., 2011. T-lex: a program for fast and accurate assessment of transposable element presence using next-generation sequencing data. *Nucleic acids research*, **39**(6):e36.
- Grimaldi, G., Skowronski, J., and Singer, M. F., 1984. Defining the beginning and end of KpnI family segments. *The EMBO journal*, **3**(8):1753–9.
- Hancks, D. C., Ewing, A. D., Chen, J. E., Tokunaga, K., and Kazazian, H. H., 2009. Exon-trapping mediated by the human retrotransposon SVA. *Genome research*, **19**(11):1983–91.
- Hancks, D. C., Goodier, J. L., Mandal, P. K., Cheung, L. E., and Kazazian, H. H., 2011. Retrotransposition of marked SVA elements by human L1s in cultured cells. *Human Molecular Genetics*, **20**(17):3386–3400.
- Harris, C. R., Dewan, A., Zupnick, A., Normart, R., Gabriel, A., Prives, C., Levine, A. J., and Hoh, J., 2009. p53 responsive elements in human retrotransposons. *Oncogene*, **28**(44):3857–65.

- Hormozdiari, F., Alkan, C., Ventura, M., Hajirasouliha, I., Malig, M., Hach, F., Yorukoglu, D., Dao, P., Bakhshi, M., Sahinalp, S. C., *et al.*, 2010a. Alu repeat discovery and characterization within human genomes. *Genome Res.*, :1–36.
- Hormozdiari, F., Hajirasouliha, I., Dao, P., Hach, F., Yorukoglu, D., Alkan, C., Eichler, E. E., and Sahinalp, S. C., 2010b. Next-generation variationhunter: combinatorial algorithms for transposon insertion discovery. *Bioinformatics*, **26**(12):i350–i357.
- Huang, C. R. L., Schneider, A. M., Lu, Y., Niranjan, T., Shen, P., Robinson, M. A., Steranka, J. P., Valle, D., Civin, C. I., Wang, T., *et al.*, 2010. Mobile interspersed repeats are major structural variants in the human genome. *Cell*, **141**(7):1171–1182.
- Iskow, R. C., McCabe, M. T., Mills, R. E., Torene, S., Pittard, W. S., Neuwald, A. F., Meir, E. G. V., Vertino, P. M., and Devine, S. E., 2010. Natural mutagenesis of human genomes by endogenous retrotransposons. *Cell*, **141**(7):1253–1261.
- Kano, H., Godoy, I., Courtney, C., Vetter, M. R., Gerton, G. L., Ostertag, E. M., and Kazazian, H. H., 2009. L1 retrotransposition occurs mainly in embryogenesis and creates somatic mosaicism. *Genes & development*, **23**(11):1303–12.
- Kazazian, H. H., Wong, C., Youssoufian, H., Scott, A. F., Phillips, D. G., and Antonarakis, S. E., 1988. Haemophilia a resulting from de novo insertion of Ll sequences represents a novel mechanism for mutation in man. *Nature*, **332**(6160):164–6.
- Khan, H., Smit, A., and Boissinot, S., 2006. Molecular evolution and tempo of amplification of human LINE-1 retrotransposons since the origin of primates. *Genome research*, **16**(1):78–87.
- Kinomoto, M., Kanno, T., Shimura, M., Ishizaka, Y., Kojima, A., Kurata, T., Sata, T., and Tokunaga, K., 2007. All APOBEC3 family proteins differentially inhibit LINE-1 retrotransposition. *Nucleic acids research*, **35**(9):2955–64.
- Kuwabara, T., Hsieh, J., Muotri, A., Yeo, G., Warashina, M., Lie, D. C., Moore, L., Nakashima, K., Asashima, M., and Gage, F. H., *et al.*, 2009. Wnt-mediated activation of NeuroD1 and retro-elements during adult neurogenesis. *Nature neuroscience*, **12**(9):1097–105.
- Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., *et al.*, 2001. Initial sequencing and analysis of the human genome. *Nature*, **409**(6822):860–921.
- Lee, E., Iskow, R., Yang, L., Gokcumen, O., Haseley, P., Luquette, L. J., Lohr, J. G., Harris, C. C., Ding, L., Wilson, R. K., *et al.*, 2012. Landscape of Somatic Retrotransposition in Human Cancers. *Science (New York, N.Y.)*, .
- Li, H. and Durbin, R., 2009. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*, **25**(14):1754–60.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and Durbin, R., 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*, **25**(16):2078–9.
- Luan, D. D., Korman, M. H., Jakubczak, J. L., and Eickbush, T. H., 1993. Reverse transcription of r2bm rna is primed by a nick at the chromosomal target site: a mechanism for non-ltr retrotransposition. *Cell*, **72**(4):595–605.

- Miki, I. N. A. H. Y. M. J. U. K. W. K. B. V. Y. N. Y., 1992. Disruption of the *apc* gene by a retrotransposal insertion of *l1* sequence in a colon cancer. *Cancer Res*, **52**:643–645.
- Muckenfuss, H., Hamdorf, M., Held, U., Perkovic, M., Löwer, J., Cichutek, K., Flory, E., Schumann, G. G., and Münk, C., 2006. APOBEC3 proteins inhibit human LINE-1 retrotransposition. *The Journal of biological chemistry*, **281**(31):22161–72.
- Muotri, A. R., Chu, V. T., Marchetto, M. C. N., Deng, W., Moran, J. V., and Gage, F. H., 2005. Somatic mosaicism in neuronal precursor cells mediated by *l1* retrotransposition. *Nature*, **435**(7044):903–10.
- Muzny, D. M., Bainbridge, M. N., Chang, K., Dinh, H. H., Drummond, J. A., Fowler, G., Kovar, C. L., Lewis, L. R., Morgan, M. B., Newsham, I. F., *et al.*, 2012. Comprehensive molecular characterization of human colon and rectal cancer. *Nature*, **487**(7407):330–337.
- Ogino, S., Galon, J., Fuchs, C. S., and Dranoff, G., 2011. Cancer immunology—analysis of host and tumor factors for personalized medicine. *Nature reviews. Clinical oncology*, **8**(12):711–9.
- Pavlíček, A., Paces, J., Zíka, R., and Hejnar, J., 2002. Length distribution of long interspersed nucleotide elements (LINEs) and processed pseudogenes of human endogenous retroviruses: implications for retrotransposition and pseudogene detection. *Gene*, **300**(1-2):189–94.
- Quinlan, A. R., Clark, R. A., Sokolova, S., Leibowitz, M. L., Zhang, Y., Hurles, M. E., Mell, J. C., and Hall, I. M., 2010. Genome-wide mapping and assembly of structural variant breakpoints in the mouse genome. *Genome Res*, **20**(5):623–635.
- Scott, A. F., Schmeckpeper, B. J., Abdelrazik, M., Comey, C. T., O’Hara, B., Rossiter, J. P., Cooley, T., Heath, P., Smith, K. D., and Margolet, L., *et al.*, 1987. Origin of the human *l1* elements: proposed progenitor genes deduced from a consensus dna sequence. *Genomics*, **1**(2):113–25.
- Selme, M. d. C., Vetter, M. R., Cordaux, R., Bastone, L., Batzer, M. A., and Kazazian, H. H., 2006. Extensive individual variation in L1 retrotransposition capability contributes to human genetic diversity. *Proceedings of the National Academy of Sciences of the United States of America*, **103**(17):6611–6.
- Stenglein, M. D. and Harris, R. S., 2006. APOBEC3B and APOBEC3F inhibit L1 retrotransposition by a DNA deamination-independent mechanism. *The Journal of biological chemistry*, **281**(25):16837–41.
- Stewart, C., Kural, D., Strömberg, M. P., Walker, J. A., Konkel, M. K., Stütz, A. M., Urban, A. E., Grubert, F., Lam, H. Y. K., Lee, W.-P., *et al.*, 2011. A comprehensive map of mobile element insertion polymorphisms in humans. *PLoS Genet*, **7**(8):e1002236.
- Taniguchi-Ikeda, M., Kobayashi, K., Kanagawa, M., Yu, C.-c., Mori, K., Oda, T., Kuga, A., Kurahashi, H., Akman, H. O., DiMauro, S., *et al.*, 2011. Pathogenic exon-trapping by SVA retrotransposon and rescue in Fukuyama muscular dystrophy. *Nature*, **478**(7367):127–31.
- TCGA Research Network, 2008. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, **455**(7216):1061–8.
- TCGA Research Network, 2011. Integrated genomic analyses of ovarian carcinoma. *Nature*, **474**(7353):609–15.

- Wang, J., Song, L., Grover, D., Azrak, S., Batzer, M. A., and Liang, P., 2006. dbrip: a highly integrated database of retrotransposon insertion polymorphisms in humans. *Hum Mutat*, **27**(4):323–9.
- Wang, T., Zeng, J., Lowe, C. B., Sellers, R. G., Salama, S. R., Yang, M., Burgess, S. M., Brachmann, R. K., and Haussler, D., 2007. Species-specific endogenous retroviruses shape the transcriptional network of the human tumor suppressor protein p53. *Proc Natl Acad Sci USA*, **104**(47):18613–8.
- Witherspoon, D. J., Xing, J., Zhang, Y., Watkins, W. S., Batzer, M. A., and Jorde, L. B., 2010. Mobile element scanning (me-scan) by targeted high-throughput sequencing. *BMC Genomics*, **11**:410.
- Woodcock, D. M., Lawler, C. B., Linsenmeyer, M. E., Doherty, J. P., and Warren, W. D., 1997. Asymmetric methylation in the hypermethylated CpG promoter region of the human L1 retrotransposon. *The Journal of biological chemistry*, **272**(12):7810–6.
- Yoder, J. A., Walsh, C. P., and Bestor, T. H., 1997. Cytosine methylation and the ecology of intragenomic parasites. *Trends in genetics : TIG*, **13**(8):335–40.
- Zerbino, D. R. and Birney, E., 2008. Velvet: algorithms for de novo short read assembly using de bruijn graphs. *Genome Res*, **18**(5):821–9.
- Zhu, J., Sanborn, J. Z., Benz, S., Szeto, C., Hsu, F., Kuhn, R. M., Karolchik, D., Archie, J., Lenburg, M. E., Esserman, L. J., *et al.*, 2009. The UCSC Cancer Genomics Browser. *Nature methods*, **6**(4):239–40.

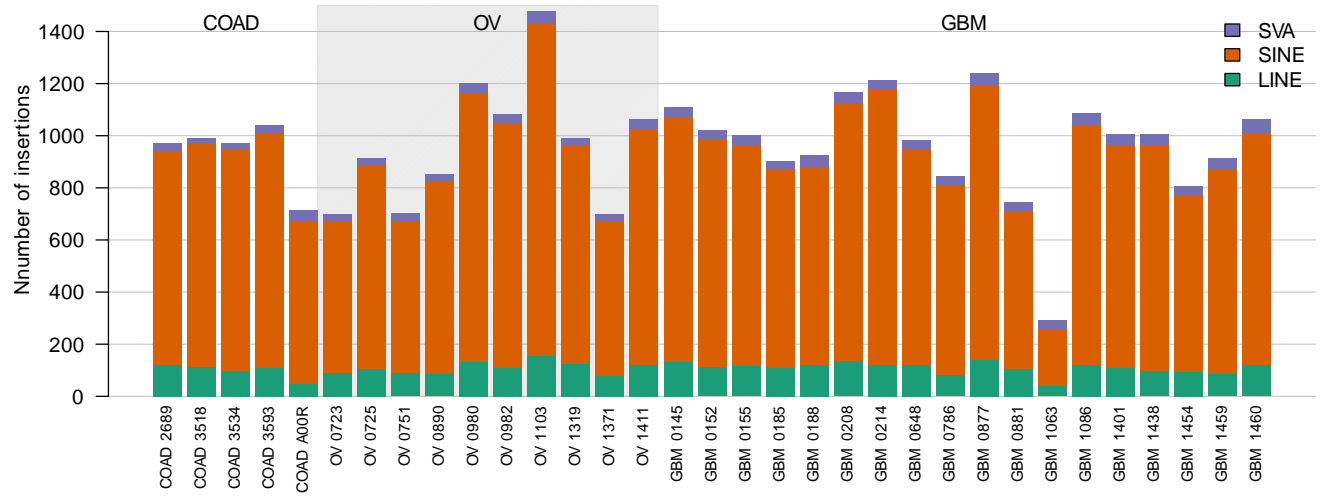


Figure 1: The number of non-reference germline insertions (found in both normal and tumor samples) per patient analyzed.

5 Figures

6 Tables

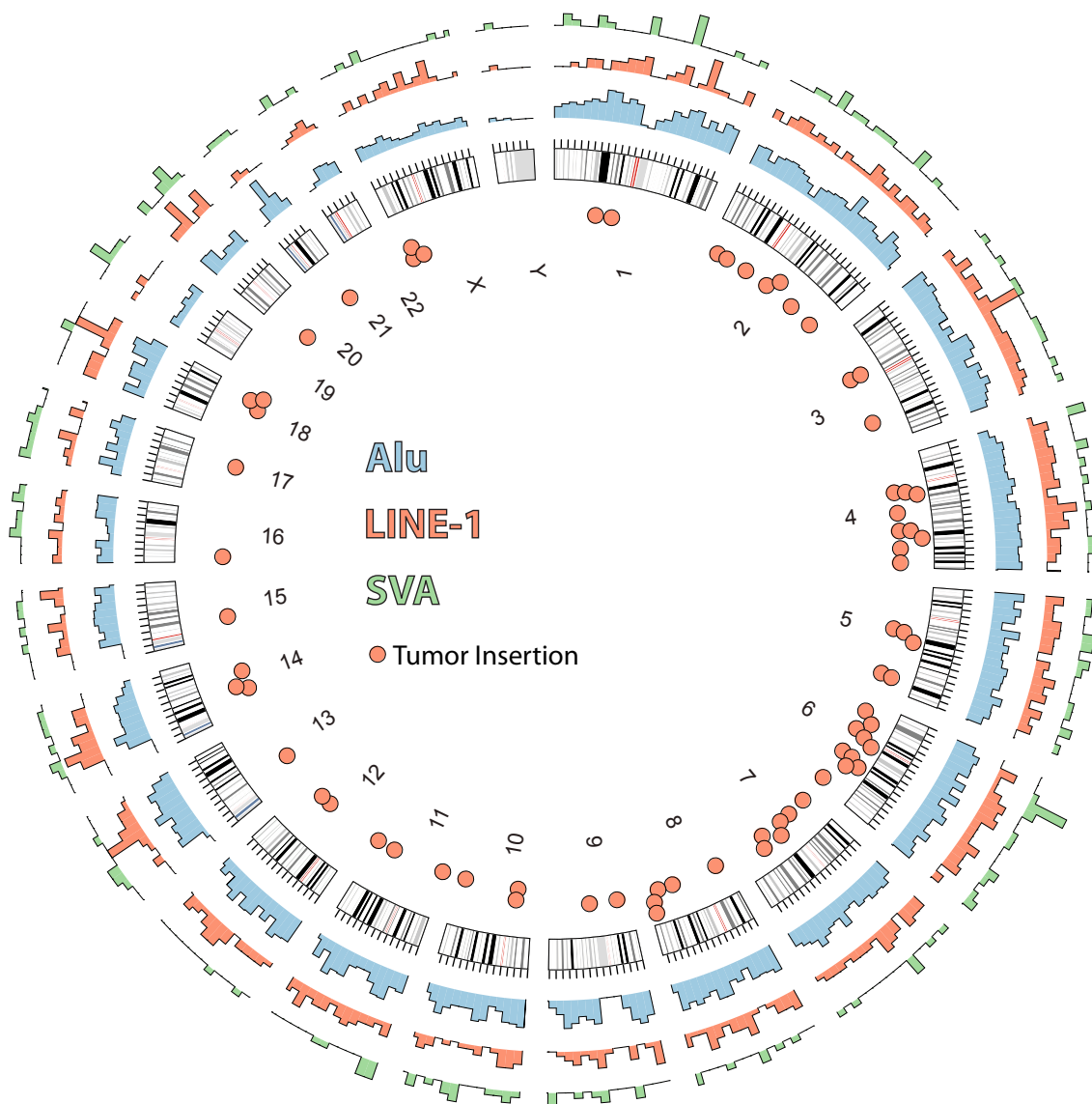


Figure 2: Plot depicting distribution of insertion site density for non-reference mobile element insertions (outer rings, colored as indicated) and tumor-specific insertions (inner ring, circles).

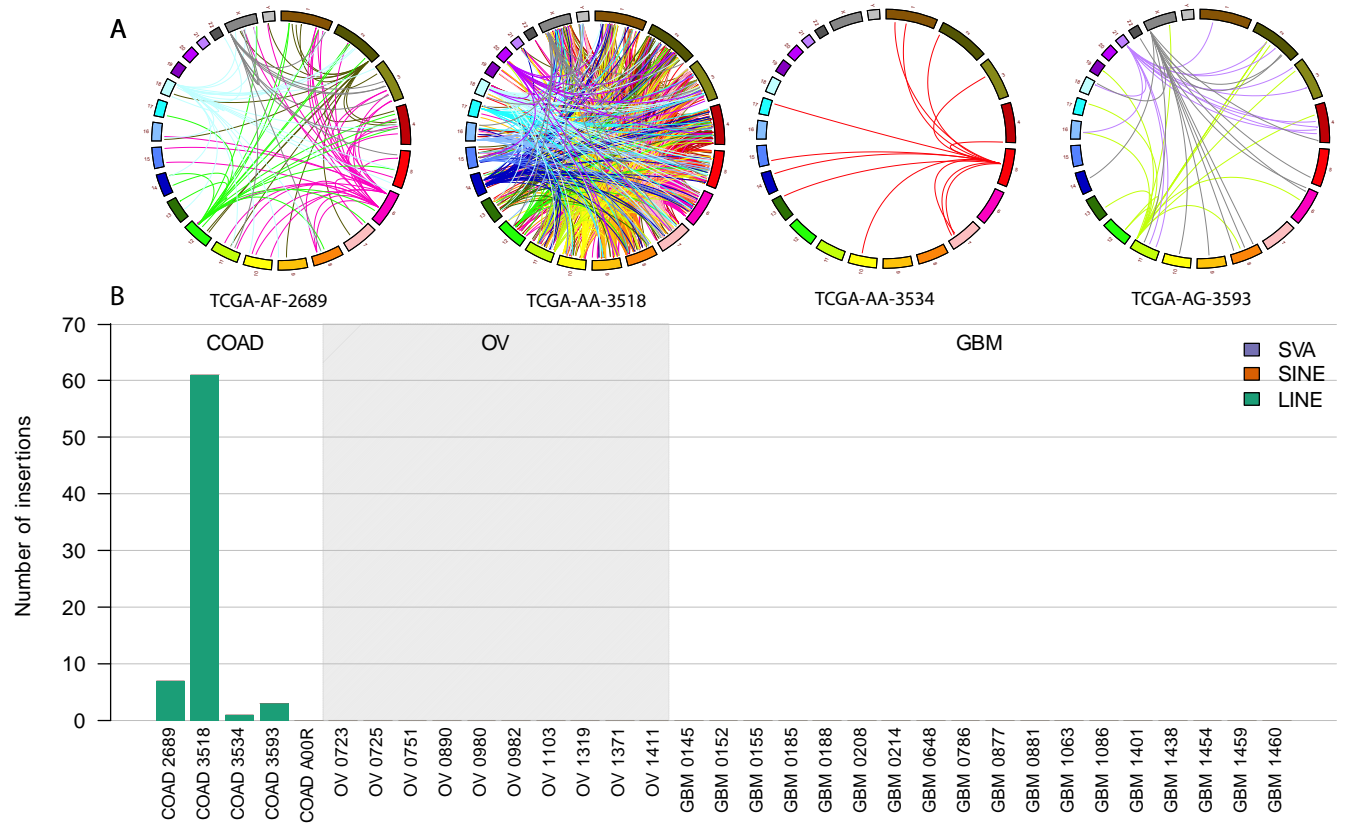


Figure 3: (A) Discordant read “one-end repeat (OER)” mappings for the 4 colorectal adenocarcinoma samples with tumor-specific retrotransposon activity. Links are shown in the color of the chromosome where the insertion occurred. (B) The number of retrotransposon insertions detected only in the tumor tissue for each of 33 patients analyzed. No cancer-specific Alu or SVA insertions were detected, and no insertions were found in ovarian carcinoma (OV) or glioblastoma multiforme (GBM), but between 1 and 61 cancer-specific insertions were found across 4 colorectal adenocarcinoma (COAD) patients tumors.

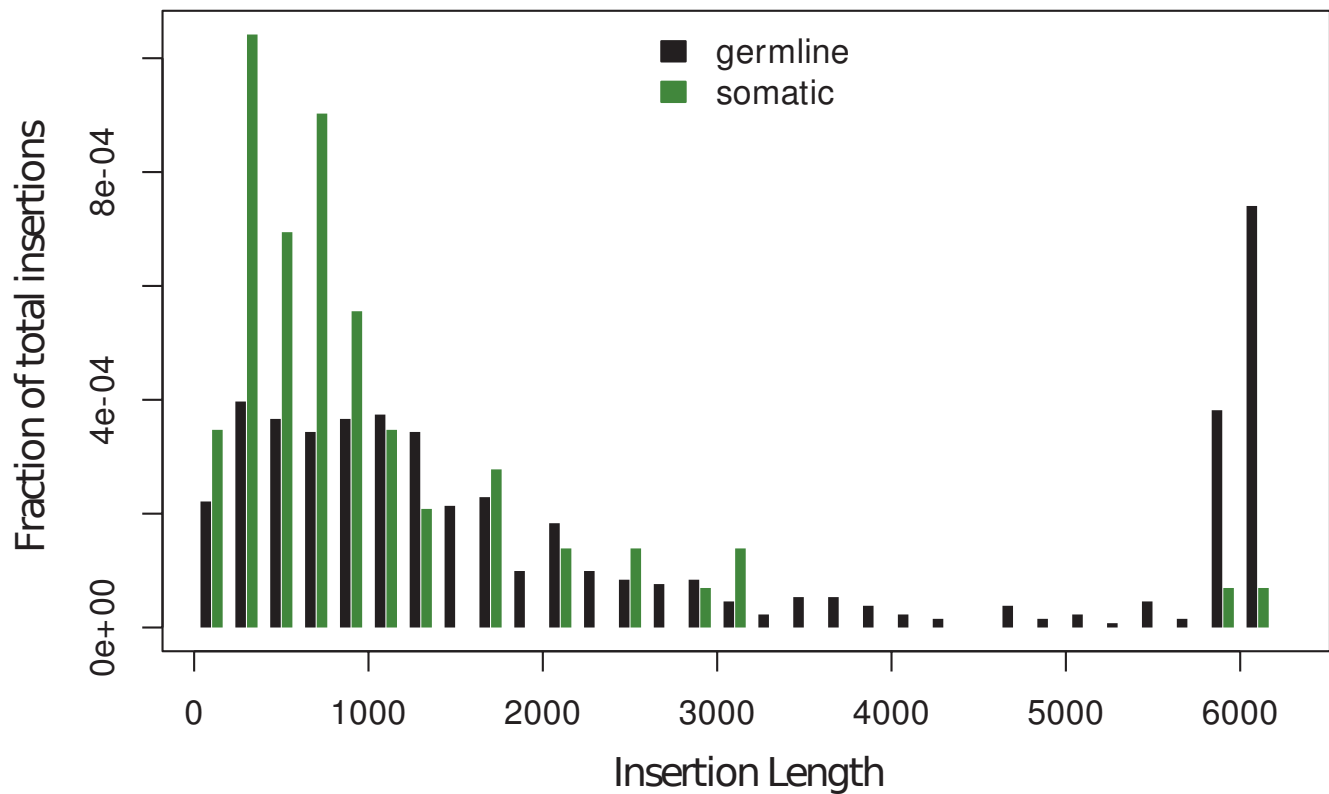


Figure 4: The lengths of non-reference L1 insertions detected in this study, binned in 200bp intervals. Germline insertions are found in the normal and the tumor tissue of a patient, found in a previous study, or found in multiple patients, and somatic/cancer insertions are found only in the tumor tissue of a single patient. There is a significant difference between the proportion of L1 germline insertions that are full length ($>5.8\text{kbp}$) and the proportion of somatic/cancer L1 insertions that are full length ($p\text{-value}=9.68 \times 10^{-6}$).

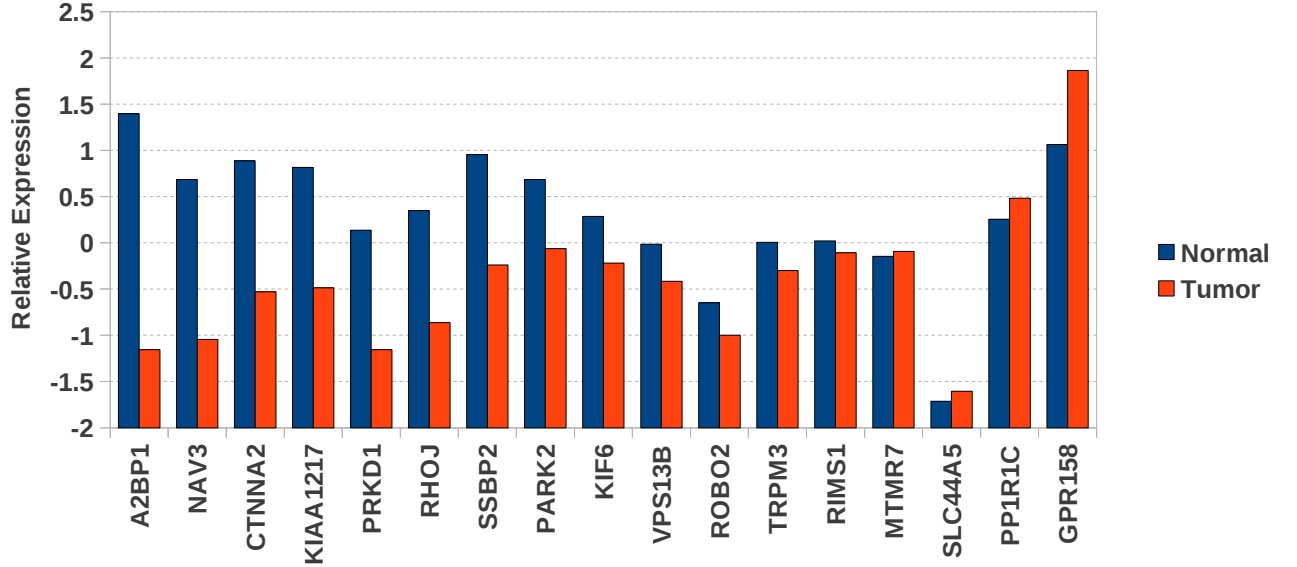


Figure 5: Expression level of genes in normal tissue (blue) prior to intronic LINE-1 insertion, and after (red). Relative expression values are taken from the UCSC Cancer Genome Browser (Zhu et al. 2009), where they were normalized by centering to the mean expression level.

LINE elements most closely related to inserted sequences

	Location	Repeat Family	Length	Best Related Insertions	Related Insertions	Characteristics
1	chr10:107127095-107133125	L1HS	6030	14	24	ORF1 broken, ORF2 intact
2	chr11:92793801-92799845	L1HS	6044	12	33	intact
3	chr11:24306074-24312123	L1HS	6049	12	28	intact
4	chr17:65966693-65972723	L1HS	6030	10	26	ORF1 intact, ORF2 broken
5	chr11:60608423-60610418	L1HS	1995	10	18	truncated
6	chr7:49690411-49696442	L1HS	6031	10	18	intact
7	chr16:22618776-22619548	L1HS	772	8	14	truncated
8	chr18:43440743-43446771	L1HS	6028	7	19	ORF1 intact, ORF2 broken
9	chr4:182113842-182114873	L1HS	1031	7	11	truncated
10	chr5:12250459-12251104	L1HS	645	7	9	truncated

Table 1: Contigs of inserted sequence were assembled for some cancer-specific transposable element (TE) insertions and aligned back to the reference genome using BLAT. Repeat masker annotated elements from the reference genome are listed according to the number of times they have the highest sequence similarity to an insertions contigs compared to all other repeat masker annotated elements (Best Related Insertions). The number of times an element has within 20 mismatches of the highest sequence similarity to an insertions contigs (in essence, is the most similar or a close second), is also listed (Related Insertions).